

A SYNCHRONIZABLE TRANSACTIONAL DATABASE METHOD AND SYSTEM

5 CROSS REFERENCES TO RELATED APPLICATIONS

This application claims the benefit of U.S. Provisional Application, Serial Number 60/224,270 filed on August 10, 2000.

10 Background of the Invention

This invention relates to the design and implementation of a database system that includes among its basic functions the ability to efficiently synchronize all data within a specified key range. A transactional database system is used for transaction processing.

15 This system generally consists of centrally storing information and allowing access to this information by users. User access typically consists of read, write, delete, and update interaction with the information on the database. These interactions are referred to as transactions.

The information on the database can be vital for many reasons. The database
20 might store banking information and its users may rely on it to accurately reflect information on thousands of accounts. The database might store medical information and its users may rely on it to properly diagnose or monitor the health of a patient.

Transactional systems may be accessed thousands of times a day and its information may be altered each time. It is often imperative to insure the security of this
25 information. These systems also face a variety of reasons to fail. Therefore, it is common for such databases to be backed up and the information made available for recovery if a failure is faced.

A common method of backup for a transactional database is synchronization. Synchronization typically involves a primary host and a secondary host. Transactions are

first processed on the primary database, the secondary database is periodically synchronized or made consistent with the primary database and is used as a backup and part of a recovery system. The role of primary and secondary may be dynamic making synchronization a multi-directional path.

- 5 Since multiple copies of the information are being kept, the copies must be kept consistent. To synchronize the databases they must both compare and transfer data.

A synchronizable database D is a set containing records of the form (key, value). The key field takes values from a totally ordered set K of keys. Any key in K occurs in D at most once.

- 10 The major operations supported by a synchronizable database as an abstract data type are insertion, deletion and retrieval of records, and a range synchronization operation. The first three are standard operations on databases with records. The last one is unique to synchronizable databases.

- The input to a range synchronization operation is an interval I of K and two
15 databases D_1 and D_2 . The operation basically tries to make the restrictions of D_1 and D_2 to I identical. In particular, it identifies three sets of keys, which are called the discrepancy sets K_1 , K_2 and K_{12} . These three sets are all subsets of the key interval I . Discrepancy set K_1 is the set of keys in D_1 which are not in D_2 , K_2 is the set of keys in D_2 which are not in D_1 , and K_{12} is the set of keys which are in both D_1 and D_2 but whose corresponding
20 records in the two databases differ in the value field.

- The operation calls different handler functions for each of these three sets. Typically, the handler functions for K_1 and K_2 would copy the missing records from one database to the other. The handler function for K_{12} would typically, for each key in the discrepancy set, compare the records in D_1 and D_2 that have the key and replace one of
25 them with the other.

 Since synchronization relies on the comparison and transfer of data, efficiency lies in the actions of comparing and transferring. As used herein, efficiency is an attribute of the cost of comparing and transferring data and the time intervals of comparing and transferring data.

operation a block mapping system redirects writes intended to overwrite a given logical block address, to a free physical block leaving the original data intact. At the same time this logical address is tentatively mapped to the physical block containing the new data written. The commit operation atomically switches all maps in memory and on disk so that these tentative mappings are made permanent. In this way bedrock provides transactional support at the lowest level of the system. Higher-level applications, including a synchronizable database system, may then use bedrock to build arbitrary data structures that inherit from below a transactional capability.

Bedrock operates within a single, preallocated file or in a raw disk environment and uses no outside logging facility. Following a system crash no recovery process is required. The data-store is instantly available in its most recently committed state. Bedrock uses a three-level mapping system that supports up to 2^{31} fixed length blocks. It represents an easy-to-administer and scalable substrate for many applications that require simple transactional support.

The bxtree layer features support for variable length records and keys, and all data is written in a portable (byte order independent) fashion. Maintained along with each data record is a fixed-length digest computed by a cryptographically strong function such as the well-known MD5 algorithm. Bxtree extends a conventional B+-tree by providing a function that, given a range of key values, efficiently computes the exclusive or (XOR) of the digest values associated with each record in the specified range. It does this by maintaining at each internal tree node, the XOR of all children. With this approach the MD5/XOR summary of an arbitrary range in the database may be computed in $O(\log n)$ time where n denotes the number of records in the database.

In the osynch protocol, two parties exchange digests of key ranges and subranges to rapidly identify regions of discrepancy. These discrepancies are then dealt with by exchanging data. The design of osynch and bxtree is aimed at minimizing the complexity of synchronization measured in terms of bits communicated, rounds of communication, and local computation. The protocol has a simple recursive structure and may be modified to adjust the tradeoff between these three objectives. The preferred embodiment

reflects a tradeoff that seems appropriate for today's computational and internet environment.

The number of rounds is asymptotically $O(\log n)$ to identify each region of discrepancy, and the local computation is $O(t \log^2 n)$ where t denotes the total number of discrepancies. The end result is that the total bits communicated is roughly proportional to the amount of differing data.

Large databases and other complex data structures maintained in non-volatile storage, e.g. on disk, often have certain important invariants to be preserved at all times for the data to make sense. In order to write data with the invariants always preserved, the updates must be made transactionally. It is desirable to have an abstraction of transactional memory, general enough to be used by many applications.

Brief Description of Drawings

Figure 1 depicts a synchronizable database.

Figure 2 depicts the software architecture of the invention.

Figure 3 represents the layout of a bedrock file.

Figure 4 represents the memory layout of the bedrock mapping mechanism.

Detailed Description of the Invention

Figure 2 depicts the software architecture of the invention. Presented is a synchronizable database as three modules: OSYNCH, 3, provides the synchronization facility, and the BXTREE and BEDROCK modules together comprise the summarizable database, 4. The BXTREE module, 5, implements an industrial strength B+-TREE, enhanced with summaries. Hashing the leaf data and storing the hashes in the leaves, and storing XOR's of all the children nodes along with each parent node produce the summaries. The BEDROCK, 6, module provides a transactional non-volatile storage layer, which allows for consistent, safe updates of the BXTREE, 5.

The bedrock is a transactional memory management system. It allows complex updates to be enacted (committed, synced) atomically (at once), over non-volatile block storage devices, such as hard drives, and to persist from one sync to the next. If an intermediate series of updates fails to commit, the previously committed state is guaranteed to be found where it was left by the previous, successfully completed commit. Any such updates in progress after a commit can be aborted forcefully, and then the previously committed state will be reactivated.

The bedrock layer consists of a disk block array and a mapping mechanism allowing for two "views" of the blocks. The ongoing updates correspond to shadow blocks, reflected in the shadow maps. The previously committed transaction persists, its blocks intact, its maps current and safely in place, governed by a single superblock. Once the updates are finished, they can be committed with a single atomic write of the new superblock, which swaps the current maps and the shadow ones. Each set of maps refers to their own blocks, so formerly shadow blocks now also become current.

Figure 3 depicts the layout of a bedrock file in non-volatile memory. Header, 10, consists of essential bedrock parameters, 12, and mapping mechanism, 14, in fixed size segments, followed by the actual blocks, 16.

Since the superblock, 20, must be written atomically, its size should be the minimal block size for any block device keeping it. An example of a small device block size is 512 bytes, corresponding to a disk sector on some systems. The preferred architectural assumptions for the operating system carrying the bedrock are as follows:

1. Atomic write of the aligned atomic size data (512 bytes). Once the write begins, it must finish, or not begin at all. Most disk controllers ensure when starting a write that it will go through even in the event of a power failure.

2. Aligned atomicity - a file write begins at an atomic write boundary. Thus, if one starts to write a file by writing out an atomic-size block, the write will indeed be atomic, and so will be all the subsequent atomic-size writes.

3. Non-volatile memory, e.g. on disk. Non-volatile memory means the data written to disk stays there until overwritten by a knowing and willing user (no latent alterations by the OS).

5

4. `fsync(2)` works. Memory-buffered files must be capable of being flushed to disk at will using a system call such as *Unix's* `fsync()`. The successful completion of the call guarantees the disk image is what the user expects it to be, and if an interruption occurs after an `fsync()`, all the data written before it will persist.

These assumptions allow one to build a layer system, capable of withstanding the most typical system failures - interruptions, such as power failures, disk sector failures, etc., including stacked failures.

A bedrock file consists of a fixed number of homogeneous, fixed-size blocks 16. It is the responsibility of an application to choose the number of blocks `n_blocks` and block size `block_size` before creating a bedrock, and then manage its memory in terms of the blocks. Blocks can be allocated, written, read, and freed. The user deals with the block addresses, which are simply numbers in the range `1 ... n_blocks`. Transactional semantics requires that certain *mapping* information, employed by the mapping mechanism described below, is stored on disk as well.

20

During a transaction, the original physical blocks 16 can't be overwritten until a commit or abort. A "shadow" physical block is written instead. The logical block addresses, given to the users, are separated from physical blocks, making a map array storing the correspondence, much like virtual memory: `map [logical] == physical` or is marked unallocated.

25

Figure 4 depicts the memory layout of the bedrock mapping mechanism. Since only those portions of the map, 26, which actually changed are updated, there's no contiguous, full-size "main" or "shadow" arrays; rather, sector-size groups of mapping entries serve as either, in accordance with a ping-pong bit flag in an upper-level mapping

structure. The superblock, 20, contains 512 bytes == 4096 bit flags, each capable of controlling an underlying mapping page. Following the convention of grouping map entries into 512 byte segments, and representing block addresses as unsigned long 4 byte entries, one has 128 entries per map segment. Should each superblock bit control (switch) a segment directly, one would end up with only $4096 * 128 = 524288$ bottom-level blocks. Since one can't enlarge the superblock, 20, without violating the atomicity condition, and the map entries and their segment sizes are likewise fixed, the invention introduces an intermediate mapping layer, which is called pages, 22. Thus, the overall mapping becomes three-level:

superblock → pages → map segments

Naturally, each page, 22, will have the atomic size of 4096 bits, controlling just as many bottom-level map segments. Each superblock, 20, bit now will govern one such page, 22. Thus, the total number of blocks switchable through the three-level mapping scheme is $4096^2 * 128 = 2^{31}$ blocks. This supremum neatly fits into a typical four-byte word, holding unsigned integer types with values as large as $2^{32} - 1$. The supremum is the function of two assumptions -- the atomic 512 byte sector size and 4 byte word, unsigned long int type representing block addresses. In the future, these architectural parameters are likely to change, as memory needs grow.

A version number, 24, is maintained, corresponding to the bedrock format, in each bedrock file, so it can be properly upgraded when the format evolves. The application must evaluate its needs ahead of time and create the bedrock of sufficient size, planning both the number of blocks and their size. Since every block overwritten during a transaction requires a shadow block, the total number of blocks overwritten must not exceed floor ($n/2$).

A single current version of the shadow mapping mechanism is kept in the memory record. All the three layers of mapping are present in a single instance each - the superblock, 20, the page array of pages, 22, and the map, 26. When an existing bedrock is

open, these structures are assembled in memory via a slalom procedure. First, the single superblock, 20, is read. Each bit of it tells which corresponding page, 22, to read -- since any of the two versions can be current while the other is shadow, they have no preferential names and are called simply 0 and 1 instances. The slalom proceeds by filling each memory slot in the pages array, 22, from either 0 or 1 page instance on disk, equaling the corresponding superblock bit. Once the single "current" page array, 22, is thus assembled, the slalom repeats for the map itself - its segments are recalled from either 0 or 1 disk instances, depending on the value of the controlling bit within the responsible page now.

Just upon the slalom assembly, the superblock, 20, pages, 22, and map, 26, faithfully reflect the current state of the superblock, 20. Any write operations, however, will reserve their "current" blocks, and update the map, tainting the map segments involved as "touched". The "shadow" map segments will go into the slots opposite of those they were read from during the slalom, and as the slalom assembly was governed by the page bits, the touched page bits must be flipped, too, in the shadow pages. Similarly to map/page bit flip, enacting the updated page will flip its bit in the superblock. The final superblock is thus prepared, ready for the final slide-in. Finalizing the commit requires that the operating system synchronizes its view of the file with that assumed by the bedrock module (for instance, on Unix, the fsync system call is provided to flush the buffers to disk). The actual slide-in calls fsync twice - immediately before and immediately after writing out the new superblock. Note that shadow pages and map segments can be written out in any order relative to each other, but the superblock must be written after them.

Figure 4 depicts the memory layout of the bedrock mapping mechanism 29
 $\text{map}[\text{logical}] == \text{physical}$. The sb, 30, and pp's, 35, are bit vectors: bit $sb[i]$ chooses pp_0 v. pp_1 segment of 4096 bits, in the range governed by i ; similarly bit $pp[j]$ chooses map_0 v. map_1 segment of 128 unsigned long physical entries (32 bits each) for the index range (of logical addresses) based on j . The touched bit vectors, 37, allow for selective writes of their "parents" upon commit or abort. The free physical, 40, and logical, 45, array lists

(the latter sharing its array with map) allow for constant-time slot and address allocation, respectively; `physical_freed_head` (pfh) protects slots freed in the current transaction from reuse until a sync/abort.

The slalom has been abstracted and parameterized for three different occasions: read, write, and abort. It works in cooperation with the touched arrays of flags 37, and is at the core of bedrock open, sync (commit), and abort operations. The touched flags 37, set naturally when allocating and overwriting blocks, provide for one to write out only those parts of the maps and pages that were actually updated. Similarly, an abort will reread only those sections of that were tinkered with and need restoration from disk. In addition, one can set all bits of all touched arrays 37 before an “open”, then simply make a slalom read, which consequently refills all the maps and the superblock 20. Also, since the touched arrays 37 accurately record all of the bit flags changed, nothing need be read but the map 29 itself when aborting.

Within the `bedrock_slalom` function, this ternary parameter is divided into two Boolean flags, reading and aborting. Then, the slalom goes as follows:

1. When reading, the superblock is read in first, since it has to govern the slalom-gathering of the pages.

2. The superblock is then traversed: when

- reading, the disk-current pages are read in accordance with the disk superblock - pages come from the current slots, those specified in the current superblock (just read). The page touched flag is reset for each page being read.

- writing, the disk-shadow pages are written in accordance with the shadow, i.e. memory superblock - pages go to the shadow slots, as opposed to the current ones recorded in the disk superblock

- aborting, the shadow, i.e. memory, superblock bits touched during the transaction being aborted are simply reversed if they were set during it in sb_touched.

As each page 22 is read, its governing sb_touched bit 37 is reset.

5

3. Similarly, the pages 22 are now traversed, and if aborting, simply restored by reversing those bits set in pp_touched 37 (ping-pong pages touched bits). The maps should be read back from the disk when aborting, as the original modified slots are overwritten in memory with the shadow ones. The rest is analogous to the superblock traversal.

4. Finally, in case of writing, the slide-in sequence, described above, is executed (fsync → write superblock → fsync).

Each write is given a logical address to write and a pointer to a memory buffer going to the bedrock under that address. The write checks first, whether the physical slot, already associated with the logical address given, was assigned in a previous transaction. If so, it has to be preserved, and a shadow slot is found for use in the course of this transaction. Preferably, shadow block protection should not be done twice for a logical address that already was written during this transaction.

20 This implementation supports a single view of the bedrock in the following sense. Suppose a new transaction begins and it performs a write operation to some bedrock address. If a read operation is issued on the same bedrock address, then the newly written data will be returned. That is, the writes performed during a transaction are reflected in subsequent reads operations in the transaction. In other words, all threads "see" only one
25 (most recent) view of the bedrock.

An extension of the bedrock is to allow multiple time-views. This embodiment allows applications to choose between the most recent view and the view of the bedrock just before the beginning of the current transaction. The former view would reflect the writes performed during the current transaction, while the latter would not.

5 This embodiment is useful when the bedrock is being used to implement a database. Such as when there is a writer thread that is responsible for performing all the update operations on the database. Thus, only the writer thread issues write operations on the bedrock. There is also a reader thread that is responsible for answering user queries on the database. Thus, the read thread only issues read operations on the bedrock. Here the writer thread would like to see the most recent view of the bedrock, while the reader thread would like to see the view just before the current transaction. This is because during a transaction while the writer is performing an update on the database, the database can be in an inconsistent state. The reader would not want to see the inconsistent state of the database. Instead the reader would like to answer queries from the slightly older, but consistent view of the database.

10
11
12
13
14
15
16
17
18
19
20 Extensions to bedrock may include parallel transactions. Parallel transactions allow more than one transaction to be active simultaneously. This can be used where two parallel threads of computation are operating on disjoint portions of the same bedrock. In this embodiment a bedrock with support for parallel transactions will be able to provide “all or nothing” guarantees to both threads.

21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
101
102
103
104
105
106
107
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025
1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079
1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133
1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187
1188
1189
1190
1191
1192
1193
1194
1195
1196
1197
1198
1199
1200
1201
1202
1203
1204
1205
1206
1207
1208
1209
1210
1211
1212
1213
1214
1215
1216
1217
1218
1219
1220
1221
1222
1223
1224
1225
1226
1227
1228
1229
1230
1231
1232
1233
1234
1235
1236
1237
1238
1239
1240
1241
1242
1243
1244
1245
1246
1247
1248
1249
1250
1251
1252
1253
1254
1255
1256
1257
1258
1259
1260
1261
1262
1263
1264
1265
1266
1267
1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295
1296
1297
1298
1299
1300
1301
1302
1303
1304
1305
1306
1307
1308
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349
1350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1360
1361
1362
1363
1364
1365
1366
1367
1368
1369
1370
1371
1372
1373
1374
1375
1376
1377
1378
1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1399
1400
1401
1402
1403
1404
1405
1406
1407
1408
1409
1410
1411
1412
1413
1414
1415
1416
1417
1418
1419
1420
1421
1422
1423
1424
1425
1426
1427
1428
1429
1430
1431
1432
1433
1434
1435
1436
1437
1438
1439
1440
1441
1442
1443
1444
1445
1446
1447
1448
1449
1450
1451
1452
1453
1454
1455
1456
1457
1458
1459
1460
1461
1462
1463
1464
1465
1466
1467
1468
1469
1470
1471
1472
1473
1474
1475
1476
1477
1478
1479
1480
1481
1482
1483
1484
1485
1486
1487
1488
1489
1490
1491
1492
1493
1494
1495
1496
1497
1498
1499
1500
1501
1502
1503
1504
1505
1506
1507
1508
1509
1510
1511
1512
1513
1514
1515
1516
1517
1518
1519
1520
1521
1522
1523
1524
1525
1526
1527
1528
1529
1530
1531
1532
1533
1534
1535
1536
1537
1538
1539
1540
1541
1542
1543
1544
1545
1546
1547
1548
1549
1550
1551
1552
1553
1554
1555
1556
1557
1558
1559
1560
1561
1562
1563
1564
1565
1566
1567
1568
1569
1570
1571
1572
1573
1574
1575
1576
1577
1578
1579
1580
1581
1582
1583
1584
1585
1586
1587
1588
1589
1590
1591
1592
1593
1594
1595
1596
1597
1598
1599
1600
1601
1602
1603
1604
1605
1606
1607
1608
1609
1610
1611
1612
1613
1614
1615
1616
1617
1618
1619
1620
1621
1622
1623
1624
1625
1626
1627
1628
1629
1630
1631
1632
1633
1634
1635
1636
1637
1638
1639
1640
1641
1642
1643
1644
1645
1646
1647
1648
1649
1650
1651
1652
1653
1654
1655
1656
1657
1658
1659
1660
1661
1662
1663
1664
1665
1666
1667
1668
1669
1670
1671
1672
1673
1674
1675
1676
1677
1678
1679
1680
1681
1682
1683
1684
1685
1686
1687
1688
1689
1690
1691
1692
1693
1694
1695
1696
1697
1698
1699
1700
1701
1702
1703
1704
1705
1706
1707
1708
1709
1710
1711
1712
1713
1714
1715
1716
1717
1718
1719
1720
1721
1722
1723
1724
1725
1726
1727
1728
1729
1730
1731
1732
1733
1734
1735
1736
1737
1738
1739
1740
1741
1742
1743
1744
1745
1746
1747
1748
1749
1750
1751
1752
1753
1754
1755
1756
1757
1758
1759
1760
1761
1762
1763
1764
1765
1766
1767
1768
1769
1770
1771
1772
1773
1774
1775
1776
1777
1778
1779
1780
1781
1782
1783
1784
1785
1786
1787
1788
1789
1790
1791
1792
1793
1794
1795
1796
1797
1798
1799
1800
1801
1802
1803
1804
1805
1806
1807
1808
1809
1810
1811
1812
1813
1814
1815
1816
1817
1818
1819
1820
1821
1822
1823
1824
1825
1826
1827
1828
1829
1830
1831
1832
1833
1834
1835
1836
1837
1838
1839
1840
1841
1842
1843
1844
1845
1846
1847
1848
1849
1850
1851
1852
1853
1854
1855
1856
1857
1858
1859
1860
1861
1862
1863
1864
1865
1866
1867
1868
1869
1870
1871
1872
1873
1874
1875
1876
1877
1878
1879
1880
1881
1882
1883
1884
1885
1886
1887
1888
1889
1890
1891
1892
1893
1894
1895
1896
1897
1898
1899
1900
1901
1902
1903
1904
1905
1906
1907
1908
1909
1910
1911
1912
1913
1914
1915
1916
1917
1918
1919
1920
1921
1922
1923
1924
1925
1926
1927
1928
1929
1930
1931
1932
1933
1934
1935
1936
1937
1938
1939
1940
1941
1942
1943
1944
1945
1946
1947
1948
1949
1950
1951
1952
1953
1954
1955
1956
1957
1958
1959
1960
1961
1962
1963
1964
1965
1966
1967
1968
1969
1970
1971
1972
1973
1974
1975
1976
1977
1978
1979
1980
1981
1982
1983
1984
1985
1986
1987
1988
1989
1990
1991
1992
1993
1994
1995
1996
1997
1998
1999
2000
2001
2002
2003
2004
2005
2006
2007
2008
2009
2010
2011
2012
2013
2014
2015
2016
2017
2018
2019
2020
2021
2022
2023
2024
2025
2026
2027
2028
2029
2030
2031
2032
2033
2034
2035
2036
2037
2038
2039
2040
2041
2042
2043
2044
2045
2046
2047
2048
2049
2050
2051
2052
2053
2054
2055
2056
2057
2058
2059
2060
2061
2062
2063
2064
2065
2066
2067
2068
2069
2070
2071
2072
2073
2074
2075
2076
2077
2078
2079
2080
2081
2082
2083
2084
2085
2086
2087
2088
2089
2090
2091
2092
2093
2094
2095
2096
2097
2098
2099
2100
2101
2102
2103
2104
2105
2106
2107
2108
2109
2110
2111
2112
2113
2114
2115
2116
2117
2118
2119
2120
2121
2122
2123
2124
2125
2126
2127
2128
2129
2130
2131
2132
2133
2134
2135
2136
2137
2138
2139
2140
2141
2142
2143
2144
2145
2146
2147
2148
2149
2150
2151
2152
2153
2154
2155
2156
2157
2158
2159
2160
2161
2162
2163
2164
2165
2166
2167
2168
2169
2170
2171
2172
2173
2174
2175
2176
2177
2178
2179
2180
2181
2182
2183
2184
2185
2186
2187
2188
2189
2190
2191
2192
2193
2194
2195
2196
2197
219

two independently usable modules - the bxtree module and the osynch module. The two modules talk to each other through a very clean interface consisting of just two functions.

5 The bxtree module implements a B+-tree based database engine augmented with functions for supporting range synchronization. For example, a typical operation on the database would input a key range (i.e. an interval of the space of keys) and return a short summary. The summary would typically contain a digest of the records of interest of all the records in the database whose keys lie in the given key range. In a bxtree, the records are stored only in the leaves, and all the leaves are at the same height in the tree. With each record a fixed size digest of the record is also stored. Each internal node also stores, for each of its children, the XOR of the digests of all the records in the subtree rooted at the child. Note that since XOR is an associative operation, this information can be efficiently maintained across tree-rebalancing operations. Because of these digests stored in the internal nodes, interval summary computation operations can be performed in time proportional to the height of the tree.

10 This embodiment uses the XOR function to combine the digests of sets of records. In fact, any "associative" function can be used instead. There is a small probability that the synchronization algorithm is unable to identify a discrepancy between the two databases. The choice of the combining function is one of the factors that determine this probability. Some functions might be more appropriate to use than others.

20 Formally, one requires an associative function that takes two p bit strings and produces another p bit string as an output. Some different embodiments for the combining function are:

addition: the function views the input and output bit strings as integers. The output is the sum of the inputs in 1's complement or 2's complement, and;

25 matrix multiplication: the function views the input and output bit strings as $m \times n$ binary matrices ($p = mn$) for some m, n . If I_1 and I_2 are the $m \times n$ matrices corresponding to the inputs, then the $m \times n$ matrix

corresponding to the output is $I_1 \times H \times I_2$ where H is a fixed $n \times m$ matrix. The function is parameterized by H .

5 The osynch (or object synchronization) module, implements the range synchronization operation on bxtree databases. This embodiment is specially tuned for the case when the databases being synchronized are located on different processors connected via a limited bandwidth link. Thus, one of the goals is to try to minimize the network traffic generated by the synchronization operation. The synchronization operation works in a number of communication rounds. In each round, the key range of interest is partitioned into smaller sub-ranges. For each sub-range, the two databases compute the summary of records lying in that sub-range and one of the databases sends its summaries to the other side. The corresponding summaries from the two sides are compared and the operation is recursively applied to sub-ranges whose summaries do not match. Only those records are transferred from one side to the other which (1) are missing on the other side, or (2) have a mismatching record on the other side. Thus, unnecessary transfer of large amounts of data is prevented.

The Bxtree Module

20 On an abstract level, the bxtree module simply implements a database engine for storing and managing records of the form (key, value), enhanced with some support for range synchronization. It provides functions for insertion, deletion and retrieval of records. In addition to these, it provides the following two functions which are used by the osynch module for implementing the range synchronization operation.

25 **Get_All_Hashes:** The input is an interval I of K . The output is a list of pairs of the form (key, hash). The list has one pair for each record in the database whose key field belongs to I . The first element in the pair is the key field of the record, and the second element is a

fixed size digest of the record. If the database has no record with key field belonging to I , an empty list is returned.

Get_Interval_Hashes: The input is an interval I of K and a positive integer H . The function partitions I into at most H disjoint sub-intervals and returns a list of triplets of the form (key_interval, num_records, hash). The list has one triplet for each sub-interval. The first element of the triplet is the sub-interval itself; the second and third elements are, respectively, the number and a fixed size digest of all the records in the database whose key fields belong to the sub-interval. Whether the database has any records with key field belonging to I or not, the list returned is always non-empty and the sub-intervals in the list form a disjoint partition of I .

This implementation of the bxtree module uses the B+-tree data structure for storing records. The internal nodes of the tree form an index over the leaves of the tree where the real data resides. In the leaf nodes where the records are stored, a fixed size digest is also stored for each record. This digest is used (1) to verify record integrity, and (2) by functions providing support for range synchronization. Each internal node stores a set of keys to guide the search for records. In addition, for each of its children, it stores a triplet of the form (address, num_records, hash) where "address" is the address of the child node, "num_records" is the number of records in the child's sub-tree, and "hash" is the XOR of the digests of the records in the child's subtree. Since XOR is an associative operation, this information can be efficiently maintained across tree-rebalancing operations.

The database allows different records to have key and value fields of different sizes. This affects the structure of the tree in several ways. It stores each node of the tree in a separate bedrock block. Since all bedrock blocks are of the same size (in bytes), each node gets the same number of bytes of storage. Thus, two full leaf nodes can have different number of records. Similarly, two full internal nodes can have different number of keys and hence different fan-out. Hence, the property (which regular B+-trees with

keys and records of fixed size exhibit) that the fan-out of any two non-root internal nodes in the tree can not differ by more than a factor of 2, is not exhibited by the tree in this implementation.

The number of nodes (or equivalently the number of bedrock blocks) accessed by an operation is a good measure of the amount of time taken. The insertion, deletion and retrieval operations make O node accesses where h is the height of the tree.

The `Get_All_Hashes` function is a simple recursive function having leaf nodes as the base case. The input to the function is an interval I in the totally ordered space of keys. A leaf node is said to be relevant for I if it contains a record with key field belonging to I , or if it is adjacent to a node like that in the left-to-right ordering of the leaf nodes. By this definition, the function only accesses the nodes that are relevant for I , and their ancestors. Thus it makes $O(h + m)$ node accesses where h is the height of the tree and m is the number of leaf nodes relevant for I . Clearly, the size of the output list is an upper bound on m .

The `Get_Interval_Hashes` function helps the `osynch` module to have the sub-intervals such that the database contains almost equal amount of data in each of the sub-intervals. The balance in the tree is used in a natural way to construct such a partition. The input to the function is a key interval I , and an upper bound H on the number of sub-intervals in the partition. The function is implemented in a recursive fashion. A typical instance of the function works with an internal tree node N' , a key interval I' , and integer bound H' . First, it identifies the set S of relevant children (i.e. children whose subtrees contain nodes relevant for I' . See above for the definition of relevance). The children of N' are sorted in a natural left-to-right order. The children in S are consecutive in this left-to-right order. The function, then, partitions S into S_1, \dots, S_n (with $n = \min\{H', |S|\}$) where each S_i consists of some children in S that are consecutive in the left-to-right order. The partition is done in such a way that each S_i has almost the same number of children from S . This partition of S naturally leads to a partition of I' into sub-intervals I_1, \dots, I_n , where I_i corresponds to S_i . The end-points of the sub-intervals come from the set of keys stored in N' , except that the left (right) end-point of I_1 (I_n) is same as the left (right) end-point of I' .

The bound H' is also partitioned into $H' = h_1 + \dots + h_n$ with no two h_i 's differing by more than one. Then for each i , working on the children in S_i (recursively or otherwise, depending on whether the children are leaf nodes or not), a partition of I_i into at most h_i sub-intervals is obtained. The output list for I is formed by concatenating the sublists from each of the S_i 's. As described earlier, N stores a triplet of the form (address, num_records, hash) for each of its children. The hash fields for some of the children in S_i are used when h_i is 1.

Again the function accesses only a portion of the sub-tree formed by leaf nodes relevant for I and their ancestors. It makes $O(h + t)$ node accesses where t is a number bounded from above by the size of the output list.

This implementation of bxtree also stores an integral field called "version" with each record. Thus, a record is really a triplet of the form (key, version, value), although conceptually the version field can be thought of as part of the value field. The version field can be used by the handler functions that are invoked by the osynch module. For example, the version field can be used to keep track of the revisions made to the value field of a record, and then the handler function for K_{12} can decide to replace the record with a smaller version field (indicating a stale value field) by the other one with a larger version (indicating a more fresh value field.)

The Osynch Module

This module implements the range synchronization operation for bxtree databases. As described earlier, the preferred implementation is designed for the case where the two databases are not located on the same processor. The two databases are assumed to be located on different processors connected via a limited bandwidth link. Thus, a major goal here is to try to minimize the traffic generated by the synchronization operation. Clearly, the problem at hand is a special case of the problem of synchronizing two bit strings located on different processors while minimizing the number of bits transferred across the network. In most situations, however, it is also desirable to minimize (or at

least keep within reasonable limits) the number of communication rounds taken by the synchronization protocol. The synchronization algorithm tries to achieve both of these objectives.

5 In this implementation, the summary of a set of records consists of the number and a fixed size digest of records in the set of portions of one database are sent across the network in order to identify the discrepancies. Once the discrepancies are identified, only those records are transferred which need to be transferred to make the databases synchronized.

10 This implementation of the synchronization operation is asymmetric in the sense that it sees one database as a local database and the other one as a remote database. It is assumed that the local database can be accessed with no (or minimal) network traffic, whereas remote database accesses generate network traffic. This asymmetric view allows the module to make optimizations that minimize network traffic. Typically, the synchronization operation will be used by a processor to synchronize its local database with some remote database.

15 The synchronization algorithm starts by asking both databases to compute a single summary of all records lying in the given key interval. The `Get_Interval_Hashes` function is invoked for this. The remote summary is transferred to the local side and compared with the local summary. If the summaries match, it is concluded that the databases are already synchronized restricted to the given key interval. Otherwise the size of the remote database restricted to the given key interval is checked. If the remote database only has a small number of records lying in the key interval, then digests for all those individual records are transferred from the remote to the local side (the `Get_All_Hashes` function is invoked here), and a record-by-record comparison is made to identify
20 discrepancies. Otherwise the remote database is asked (by calling the
25 `Get_Interval_Hashes` function) to partition the key range into smaller sub-intervals and send summaries for each of the sub-intervals. These remote summaries are then compared against corresponding local summaries and the operation is invoked recursively for sub-interval whose summaries do not match.

5 The synchronization operation takes at most $O(\log n)$ communication rounds to identify each discrepancy, where n is the total size of the two databases restricted to the given key interval. Thus, the total number of communication rounds is $O(t \log n)$, where t is the combined size of the three discrepancy sets. Also, since all btree operations take time proportional to the height of the tree, the overall computational burden for the synchronization operation is $O(t \log^2 n)$.

10
15
20
25
30
35
40
45
50
55
60
65
70
75
80
85
90
95
100
105
110
115
120
125
130
135
140
145
150
155
160
165
170
175
180
185
190
195
200
205
210
215
220
225
230
235
240
245
250
255
260
265
270
275
280
285
290
295
300
305
310
315
320
325
330
335
340
345
350
355
360
365
370
375
380
385
390
395
400
405
410
415
420
425
430
435
440
445
450
455
460
465
470
475
480
485
490
495
500
505
510
515
520
525
530
535
540
545
550
555
560
565
570
575
580
585
590
595
600
605
610
615
620
625
630
635
640
645
650
655
660
665
670
675
680
685
690
695
700
705
710
715
720
725
730
735
740
745
750
755
760
765
770
775
780
785
790
795
800
805
810
815
820
825
830
835
840
845
850
855
860
865
870
875
880
885
890
895
900
905
910
915
920
925
930
935
940
945
950
955
960
965
970
975
980
985
990
995
1000

In practice however, the number of rounds taken and the network traffic generated depend on several factors including how the discrepancies are distributed across the entire key range and how the parameters in the algorithm are chosen. There are two main parameters in the synchronization algorithm which need to be carefully chosen. The first one determines when to invoke the Get_All_Hashes operation instead of further partitioning with the Get_Interval_Hashes operation. The second parameter determines the number of partitions obtained through calls to Get_Interval_Hashes. The choice of these parameters to a very large extent determines the actual network traffic generated and the number of communication rounds taken. For example, if the algorithm decides to invoke the Get_All_Hashes function on key intervals for which the remote database has a huge number of records, then the number of communication rounds would be small but every round would generate heavy network traffic. Similarly, if the Get_All_Hashes function is invoked only on very small key intervals, then the number of rounds will be large. Note that the objective of the synchronization protocol is not only to minimize the amount of network traffic generated, but also to keep the number of communication rounds within reasonable limits. These two objectives often compete with each other. Thus, the two parameters above need to be tuned to suit different applications and conditions.

25 Associative functions other than XOR may be used to reduce any error possibilities. The error probability can be further reduced by using more than one digest. In other words, the summary used by the synchronization algorithm could contain a set of digests as opposed to just one digest.

